

Building data.eurecom.fr

Anne-Elisabeth Gazet

Fall 2011

Introduction

Context

The web as we know it today mainly consists of documents that are linked together. We are now moving to a new era, where raw data is being published on the web as *Linked Data*, and woven into the *Web of Data* or *Semantic Web*.

“The vision of the Semantic Web is to extend principles of the Web from documents to data. Data should be accessed using the general Web architecture using, e.g., URI-s; data should be related to one another just as documents (or portions of documents) are already. This also means creation of a common framework that allows data to be shared and reused across application, enterprise, and community boundaries, to be processed automatically by tools as well as manually, including revealing possible new relationships among pieces of data.”

from the W3C Semantic Web FAQ¹

In order to achieve these goals, new technologies were developed:

- RDF, for Resource Description Framework, is the main building block of the Semantic Web. It's a simple interchange format, where data is represented as triples of the form (subject, predicate, object) ;
- OWL (Web Ontology Language) and RDFS (RDF Schema) are vocabularies for describing RDF properties and classes ;
- SPARQL (SPARQL Protocol and RDF Query Language) is a protocol and query language for semantic web data sources ;
- SPARQL Update, a companion language for SPARQL, enables modifying a data source.

More information on the Linked Data paradigm can be found in [1].

¹<http://www.w3.org/2001/sw/SW-FAQ>

Motivation for the project

There is a trend to increase government transparency by releasing more and more public sector information on the web as linked data. The first initiative of the kind was data.gov.uk in the United Kingdom, which was followed by similar projects around the world.

Universities and schools also generate a lot of data about students, promotions, professors, courses, publications, departments, rooms, schedules, exams, etc. The goal of this project is to take all this data generated by EURECOM, transform it in semantic formats (RDF), interlink it with other data (from other universities and datasets) and publish the whole as linked data in order to develop showcase applications that provide useful services to students and professors.

Contents

Chapter 1

Similar projects

This project was inspired by recent similar initiatives in other universities. I studied them in order to see what kind of data they made available, which approach they took in order to publish their data, and what kind of applications were built thanks to the newly exposed data.

data.open.ac.uk - open linked data from The Open University

The OU's LUCERO (Linking University Content for Education and Research Online) project¹ was the first initiative to expose public information from a university as Linked Open Data. The data.open.ac.uk platform was developed as part of the LUCERO project.

Datasets exposed on this platform include:

Open Research Online: publications from OU researchers ;

OU podcasts: collection of Audio and Video material related to education and research at the Open University ;

Course Descriptions ;

OpenLearn: metadata related to units of teaching and learning material openly available from the OpenLearn website ;

KMi Planet Stories: from the online news system of the Knowledge Media Institute ;

KMi People Profiles.

Vocabularies used to represent the data include:

¹<http://lucero-project.info/lb/>

BibO: the Bibliographic Ontology² is used to represent information about publications originating from OU researchers ;

W3C Media Ontology: this ontology³ is used to describe audio and video material in the OU podcasts dataset ;

AIISO and the courseware ontology: the Academic Institution Internal Structure Ontology⁴ and the courseware ontology⁵ are used to describe courses ;

FOAF: the Friend Of A Friend ontology⁶ is used to represent information on the staff members at the Knowledge Media Institute

Applications

- OpenLearn Linked Data⁷ makes use of data from data.open.ac.uk to suggest courses, podcasts and other OpenLearn units that relate to an OpenLearn Unit ;
- The OU Expert Search⁸ system allows users to find academics at the Open University who are experts in a given domain ;
- OUExperts⁹ is a mobile (android) application to find Open University experts in a given domain, and connect to their social network ;
- Buddy Study¹⁰ suggests potential contacts and Open University courses to follow for students, based on the analysis of the topics in the user's Facebook page.

²<http://bibliontology.com/specification>

³<http://www.w3.org/TR/mediaont-10/>

⁴<http://vocab.org/aiiso/>

⁵<http://courseware.rkbexplorer.com/ontologies/courseware>

⁶<http://xmlns.com/foaf/spec/>

⁷<http://fouad.zablith.org/apps/openlearnlinkeddata/>

⁸<http://kmi-web15.open.ac.uk:8080/ExpertSearchClient/> (accessible inside the OU network only)

⁹<http://vimeo.com/19743762>

¹⁰<http://www.matthew-rowe.com/BuddyStudy/>

data.southampton.ac.uk

data.ox.ac.uk

LODUM - Linked Open Data University of Münster

National Research Council (CNR) - data.cnr.it

Conclusion

Chapter 2

Data model and design of our ontology

Publishing data in RDF format does not make it automatically useful and reusable. It is very important to choose carefully the way the data will be modeled, and to document this. This is the most time-consuming task in such a project, but it is well worth it: a well-designed data model will help future application developers consume the data, and combine it with external datasets.

2.1 Available data

One of the difficulties with this project is that we had to think about the data model *before* actually getting access to the data, and therefore, before knowing exactly what data would be exposed. To have a clearer idea of what data was stored at EURECOM, I had a look at the internet and intranet sites of the institute. The information available there is mainly about:

- people (teachers, researchers, doctoral students, staff members),
- research outputs,
- courses.

Assuming these would be the datasets we would access later on, we designed an RDF data model.

2.2 Choice of external vocabularies

To achieve re-usability, it is paramount to represent data using vocabularies that are used by other projects as well. This way, the data will be easier to consume for existing applications. Furthermore, the fact that a vocabulary

was accepted by the community as a de facto standard is a good hint that this vocabulary is well designed, and well fitted to its domain.

Yet, since there are still relatively few universities publishing their data as Linked Data, we could not rely only on popularity to make choices. Besides, each dataset being unique, there are relationships that we wanted to represent, to which we did not find any equivalent in the other projects' datasets. In this case, we used tools such as Schema-Cache¹, Schemapedia² and Sindice³ to see if there existed any term we could use.

Alongside popularity, the other criteria we took into account to select vocabularies are:

- the similitude between the concepts and terms present in the vocabulary, and the concepts needed to represent the data ;
- whether the vocabulary was created as part of a bigger project: this means it is the result of a consensus rather than the vision of a sole individual on the domain.

In the end, we selected the following vocabularies:

- FOAF⁴: the Friend Of A Friend vocabulary defines terms for describing persons and their relations to other people and objects. We use it to describe people at EURECOM.
- Dublin Core terms⁵: this vocabulary is a set of generic metadata terms whose purpose is to describe numeric and physical resources. We use for example the predicate dc:creator to represent the relationship between a document and its creator.
- Participation⁶: The participation ontology is a simple model for describing the roles that people play within groups. As discussed later, we subclass the class part:Role in the REVE ontology, to describe roles played by people inside EURECOM.
- AIISO⁷: The Academic Institution Internal Structure Ontology provides classes and properties to describe the internal organizational structure of an academic institution. We use the aiiso:Course and aiiso:KnowledgeGrouping classes to define courses and tracks.

¹<http://schemacache.com/>

²<http://schemapedia.com/>

³<http://sindice.com/>

⁴<http://xmlns.com/foaf/spec/>

⁵<http://dublincore.org/documents/dcmi-terms/>

⁶<http://vocab.org/participation/schema>

⁷<http://vocab.org/aiiso/schema>

- BIBO⁸: The Bibliographic Ontology describe bibliographic things on the semantic Web in RDF. It has been inspired by many existing document description metadata formats. We use it to describe the articles in EURECOM’s scientific publications repository.
- LOD⁹: The ontology for Linking Open Descriptions of Events is an ontology for publishing descriptions of historical events as Linked Data, and is designed to be compatible with other event-related vocabularies and ontologies. We use the predicates lode:atTime, lode:atPlace and lode:involvedAgent to describe course sessions.
- OWL-Time¹⁰: This ontology provides a vocabulary for expressing facts about topological relations among instants and intervals, together with information about durations, and about datetime information. We use it to describe the temporal aspects of course sessions.
- Rooms¹¹: It’s simple vocabulary for describing the rooms in a building, which we use to describe the rooms and buildings where courses take place.

2.3 REVE - the Research and Education Vocabulary for EURECOM

Some of the concepts we need to represent the data are specific to EURECOM, so we defined a set of new terms. When it was possible, we defined these terms as extensions of existing terms.

The vocabulary is specified using the RDFS and OWL languages. The latest specification can be found in RDF format at the following address: <https://github.com/aegazet/REVE-Ontology>

2.3.1 Modeling issue: what is a course ?

One of the main modeling issues we encountered was due to a particularity of the french language: the word “cours” can mean many different things, even when we restrict ourselves to the context of a school. For instance, it can designate a lecture, or a course. We also had to make the distinction between a course which took place in a given semester, and the general subject of the course. E.g., “WebSem - Spring 2010” and “WebSem - Spring 2011” are two modules which both treated the same subject — an introduction to the Semantic Web — but they took place in different times, and different people were registered for them.

⁸<http://bibliontology.com/specification>

⁹<http://linkedevents.org/ontology/>

¹⁰<http://www.w3.org/TR/owl-time/>

¹¹<http://vocab.deri.ie/rooms>

In such cases, it is essential to discuss the issue with experts on the domain we try to model. Here we asked Alexia Cepero, the Student's Pedagogy Officer of EURECOM, what was the correct representation. This led us to define the two following classes:

- Course: *a teaching unit. A course is composed of several course sessions.*
- Course session: *a course session is a punctual event on which teacher and students gather, for a given course.*

Counter intuitive as it might seem, there is no concept representing the unity of subject between, for instance, “WebSem - Spring 2010” and “WebSem - Spring 2011”. But in the data we publish, we can still hint at a relationship between the two instances, for example by stating they have the same `aiiso:code` attribute.

2.3.2 Representing an n-ary relation

Another issue we encountered, but was easier to solve, resulted from an assumption I made on the data which turned out to be false. I had assumed the credits one can earn by successfully completing a course only depended on the course itself. It turns out that it also depends on the track of study the student is following at EURECOM.

In other words, what I had modeled as a (course, credit) relationship is actually a (course, credit, track) relationship. The minor trouble with this, is that such a relationship cannot be represented with a single triple.

This is a common situation in ontology modeling. Good modeling practices have therefore been identified by the community: ontology patterns for representing n-ary relations in RDF and OWL are presented in a W3C working group note¹². We picked the first pattern described in this note, which consists in creating a new class and n new properties to represent the n-ary relation.

2.4 First data model

This first model was developed having in mind the data I had observed on the various sites of the school. Therefore, it is much more complete than what has been implemented so far, mainly because a lot of data has not been exposed yet. This preliminary model is documented here in the hope that it will be useful for future developments of the project.

You will find as annexes several graphs illustrating this model. Each graph is focused on a particular domain: persons, documents, and courses.

¹²<http://www.w3.org/TR/swbp-n-aryRelations/>

2.4.1 Persons

This model was inspired by data that is visible on EURECOM’s public staff directory.

2.4.2 Documents

This model is mainly a subset of the Bibliographic Ontology.

Representing a list of authors

There are two options for representing an article’s list of authors through the predicate `bibo:authorList`: the range of this property is the union of the classes `rdf:Seq` and `rdf:List`. Both of these classes are *container* classes, they were designed to represent sets of objects, which is not an easy task with triples. But they are quite different:

- `rdf:List` is a linked list: instances of this class have an `rdf:first` attribute pointing at the *head* of the list, and an `rdf:rest` attribute pointing at the *tail* of the list (another `rdf:List`). The end of the list is indicated by setting the `rdf:rest` to `rdf:nil`.
- `rdf:Seq` is more like a table: an instance of this class has attributes of the form `rdf:_1`, `rdf:_2`, ... `rdf:_n` pointing at the elements of the list.

The two approaches have upsides and downsides. Using `rdf:List`, you can indicate exactly what the elements of the list are, which you cannot with `rdf:Seq`. Indeed the Semantic Web is built on an *open world assumption*: the absence of a statement does not mean that the statement is false. In other words, you can state that an `rdf:Seq` has *at least* elements, but you cannot state that there are no other elements.

On the other hand, if you use `rdf:List` to represent your data, you will need a lot of instances of `rdf:List`. This makes the structure of the data complicated, and in general hard to query. The more so as, since it makes no sense to give a globally unique identifier to every `rdf:List` in your data, this approach leads to creating *blank nodes* for them. Blank nodes are nodes without a URI, and it is good practice to avoid them whenever possible, as stated in [1]:

“The scope of blank nodes is limited to the document in which they appear, meaning it is not possible to create RDF links to them from external documents, reducing the potential for interlinking between different Linked Data sources. In addition, it becomes much more difficult to merge data from different sources when blank nodes are used, as there is no URI to serve as a common key.”

2.5 Present state of the data

Chapter 3

Data conversion

3.1 Accessing the data

3.2 Building RDF triples with Python

3.3 Enriching the data with an external source:
GeoNames

3.4 Loading triples in a triplestore using SPARQL
Update

Chapter 4

Showcase application : a map of EURECOM publications

Chapter 5

Future work

Bibliography

- [1] Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space*. Morgan & Claypool, 1st edition, 2011.